# Basic BioGrid

Rick Stevens

Argonne National Laboratory

University of Chicago
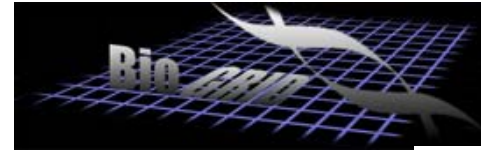
# Outline

- Biology: The Science for the 21st Century
- Sizing up the opportunities for BioGrid
- Biology is a different way of thinking
- Systems biology and whole cell modeling
- Requirements for the BioGrid
- A modest proposal
- Some recommended reading
- Conclusions

# Many BioGrid Projects

- EUROGRID BioGRID
- Asia Pacific BioGRID
- NC BioGrid
- Bioinformatics Research Network
- Osaka University Biogrid
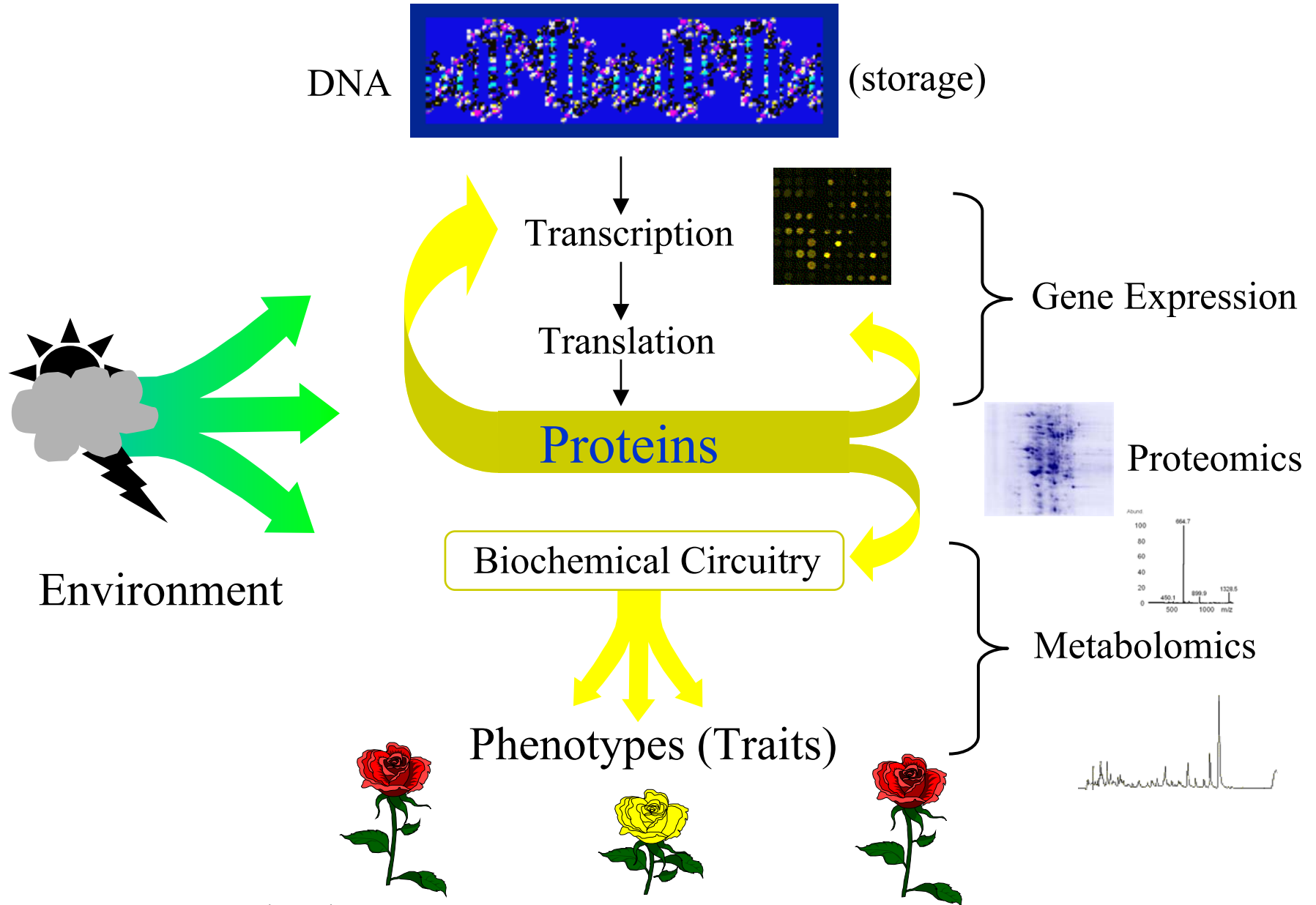- Indiana University BioArchive BioGrid

Rick Stevens

Argonne ✪ Chicago

# The New Biology

- Genomics

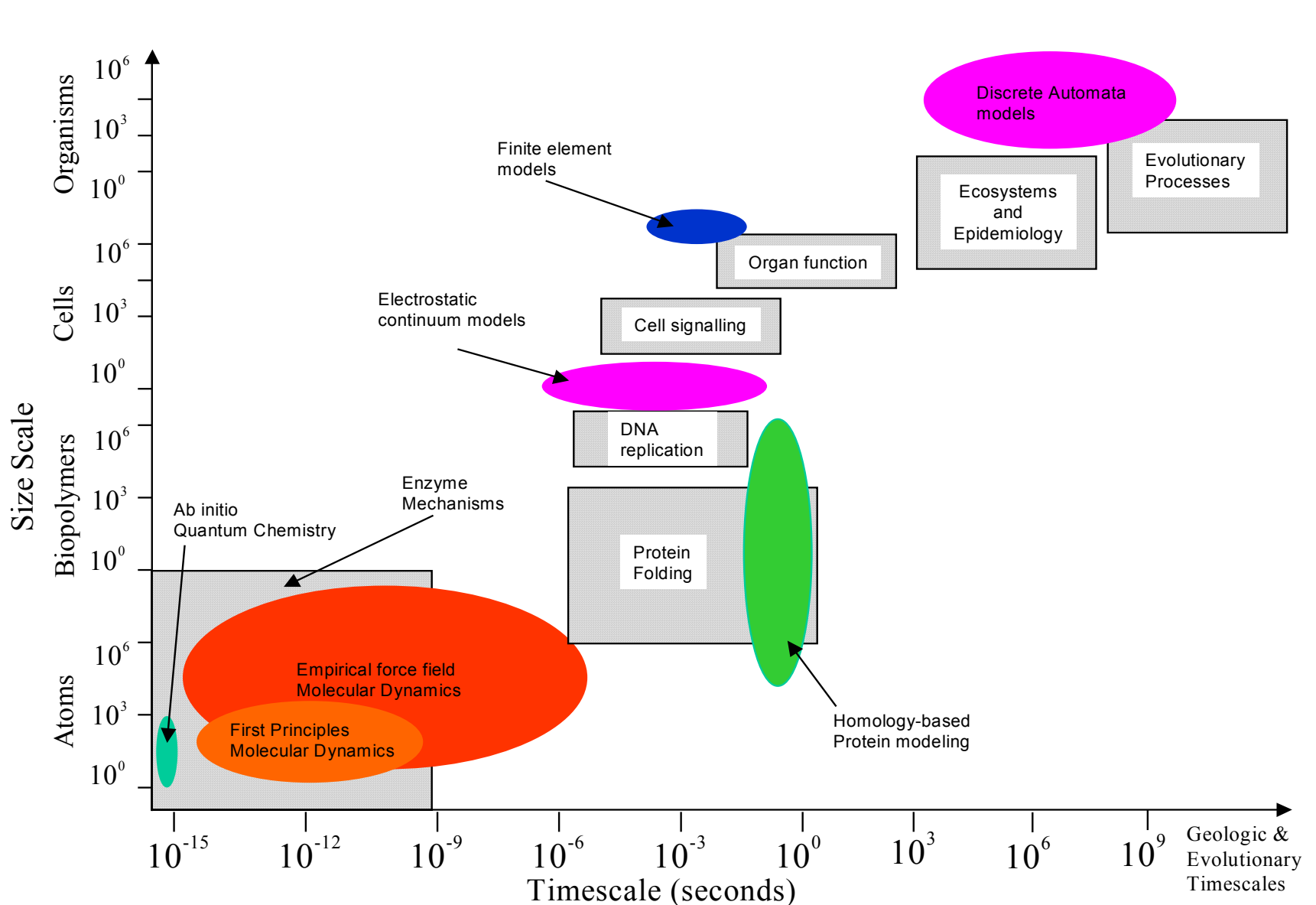- Functional Genomics

- Proteomics

- Structural Biology

- Gene Expression

- Metabolomics

- Advanced Imaging

- High-throughput methods
  - Low cost
  - Robotics

- Bioinformatics driven

- Quantitative

- Enables a systems view

- Basis for integrative understanding
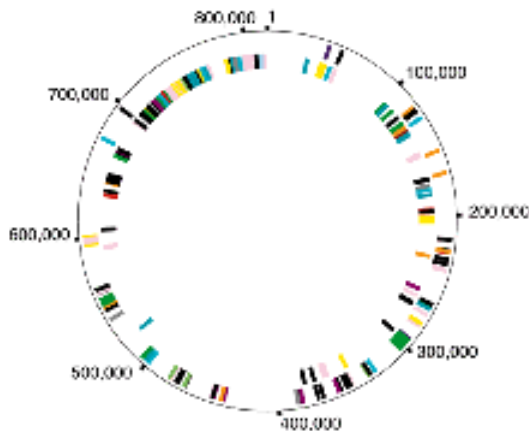  - Global state
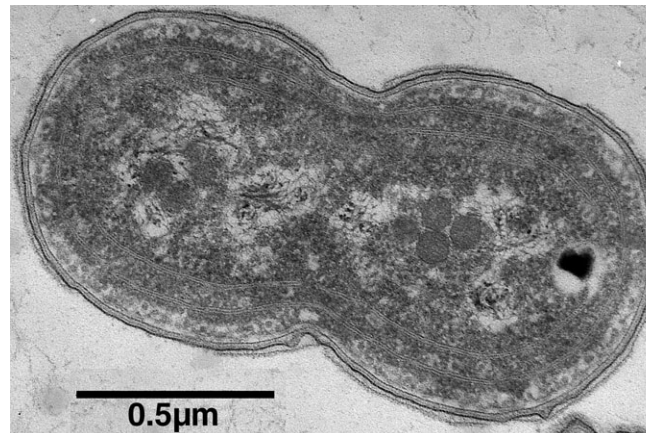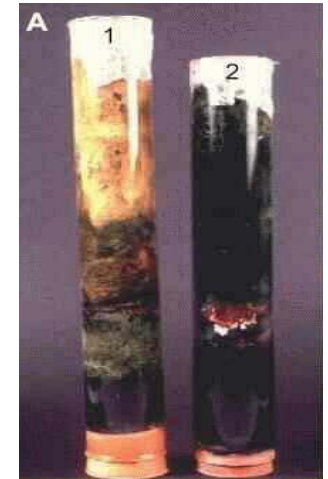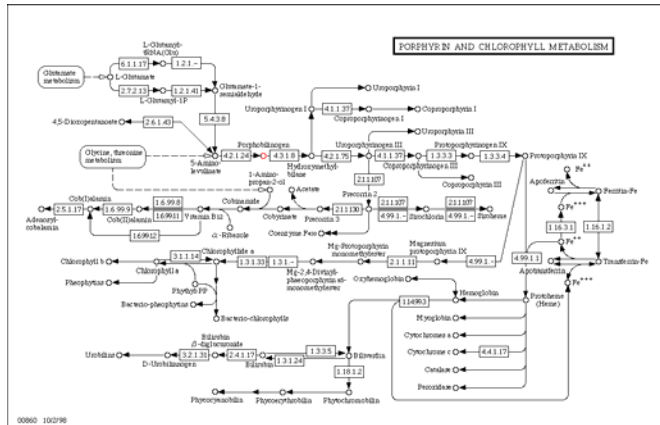  - Time dependent

# Predicting Life Processes: Reverse Engineering Living Systems



From Bruno Sobral VBI

# 24 Orders Magnitude of Spatial and Temporal Range

# Genes → Cell Networks → Organisms → Populations → Ecosystems



Rick Stevens

# Computational analysis and simulation have important roles in the study of each step in the hierarchy of biological function

**DNA sequence**

**Protein sequence and regulation**

**Protein structure**

**Protein/enzyme function**

**Sequence Annotation**

**Homology based protein structure prediction**

**Molecular simulations**

Promoter

T
A
T
A
C
A Q
G
Message
T
A Y
C
C
G R
T

**Expt. data integration**

**Organism simulations**

**Pathway simulations**

**Network analysis**

G6P

P$_i$          P$_i$

PHI

F26BPase          F16BPase

F26BP          F6P          F16BP

PFK2          PFK1

⊖          ⊖   ⊕

ADP   ATP          ATP   ADP

Citrate          AMP

**Bacterial communities & multicellular organisms**

**Bacteria and cells**

**Metabolic pathways & regulatory networks**

**Multi-protein machines**

# The New Biology in Action

**Hypothetical example:**

New data showing that a gene forms complex with DNA repair enzyme



**Fill in component in DNA repair pathway**
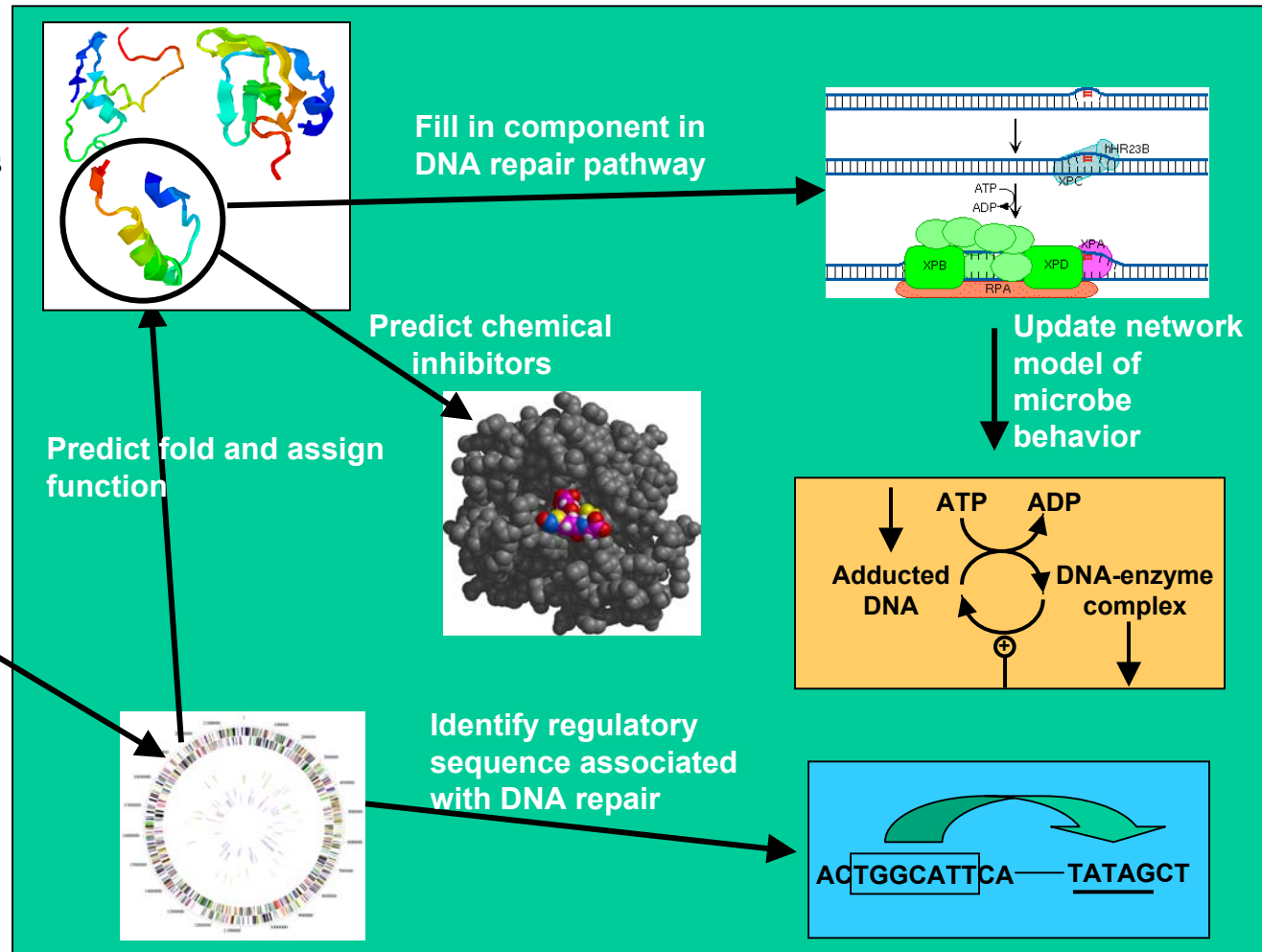
**Predict chemical inhibitors**

**Predict fold and assign function**

**Add function to genome database**

**Update network model of microbe behavior**

ATP    ADP

**Adducted DNA**    **DNA-enzyme complex**

**Identify regulatory sequence associated with DNA repair**

ACTGGCATTCA — TATAGCT

From Mike Colvin, LLNL

# An Integrated View of Simulation, Experiment, and Bioinformatics



**Problem Specification** → **Simulation** → **Browsing & Visualization**

**SIMS***

**Analysis Tools**

**Database**

**LIMS**

**Experimental Design** → **Experiment** → **Browsing & Visualization**

*Simulation Information Management System

Rick Stevens

Argonne ✶ Chicago

# Genomics is Powering the New Biology, but Computing is in the Drivers Seat

Ecological Processes and Populations

Tissue and Organismal Physiology

Cellular & Developmental Processes

Biochemical Pathways & Processes

Genes, Proteins, RNAs, and other Biomolecules

Genome Sequences

**Computation**

**Experiments**

**Large-scale Genome Sequencing**

Morphogenesis and Development

Simulation of Metabolic and Signal Transduction Pathways

Predicting Catalysis, Molecular Dynami

Structures of Multi-molecular complexes

Predicting Effects of Variation

Predicting Three-Dimensional Structures of Proteins and RNAs

Predicting Functions

Predicting Protein Sequence

Simulating and Understanding Gene Expression Networks

Genes and Gene Structures

Reconstructing Phylogeny, Homology, and Comparitive Approaches

Assembled Genomes

Sequence Variation of Populations

.cago

# Systems Biology

- Integrative (synthetic) understanding of a biological system
  - Cell, organism, community and ecosystem
- Counterpoint to reductionism
  - Requires synthesizing knowledge from multiple levels of the system
- Discovery oriented not necessarily hypothesis driven
  - Data mining vs theorem proving

# Biology is BIG!!

- 3500 Millions years of evolution
  - ~10M extant species (distinct genomes)
    - 200 genomes sequenced so far many to go!!
    - >100M extinct species
  - $10^8$ average genome size (coding region)
  - ~$10^8$ bp x ~$10^7$ sp = $10^{15}$ bp of genetic diversity
  - $10^{11}$-$10^{13}$ total genes and gene products
    - ~2,000 protein structures determined
    - typical bacteria has about 3,000 types of proteins

# FIVE KINGDOMS

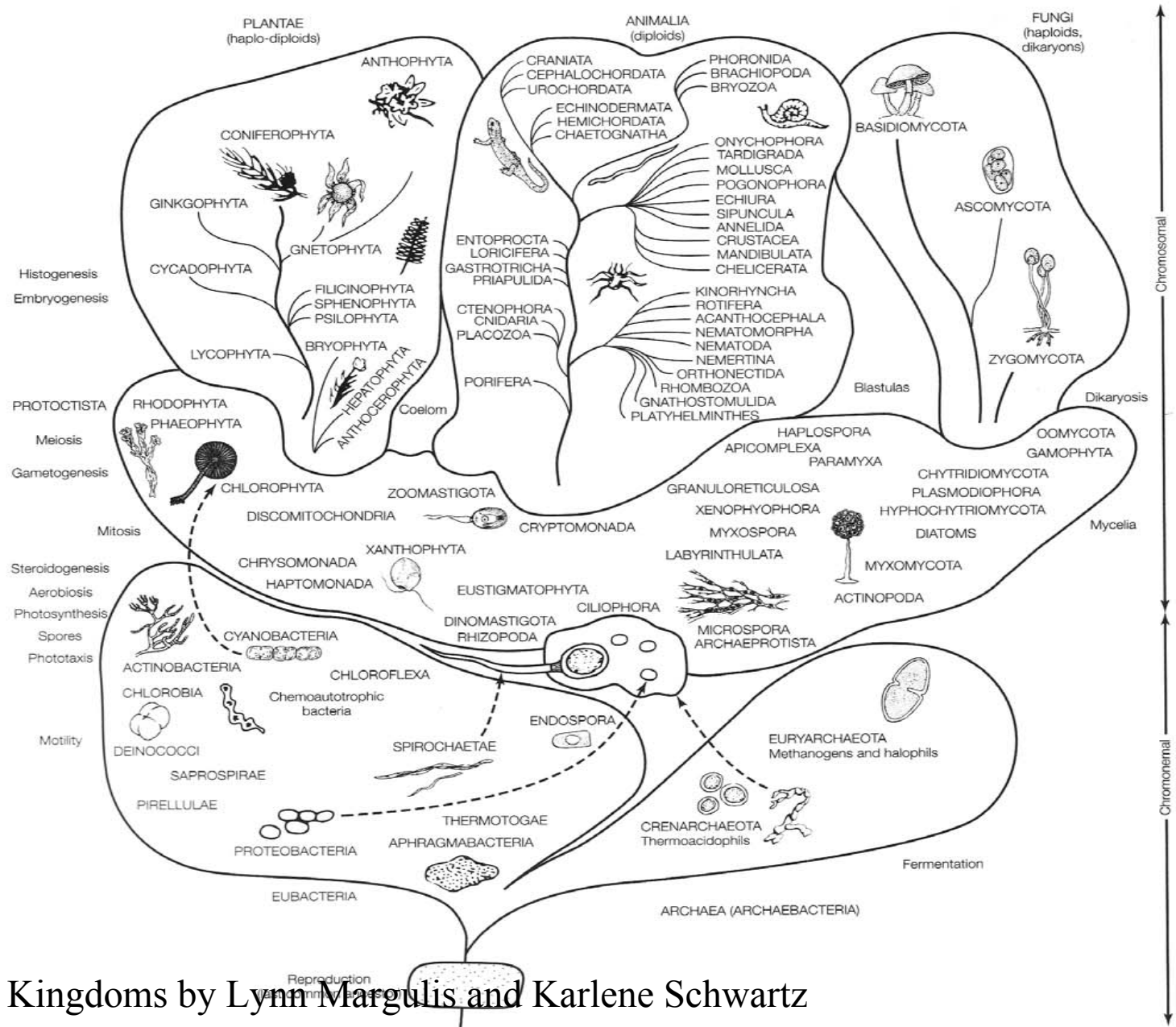An Illustrated Guide to the Phyla of Life on Earth

THIRD EDITION

LYNN MARGULIS and KARLENE V. SCHWARTZ
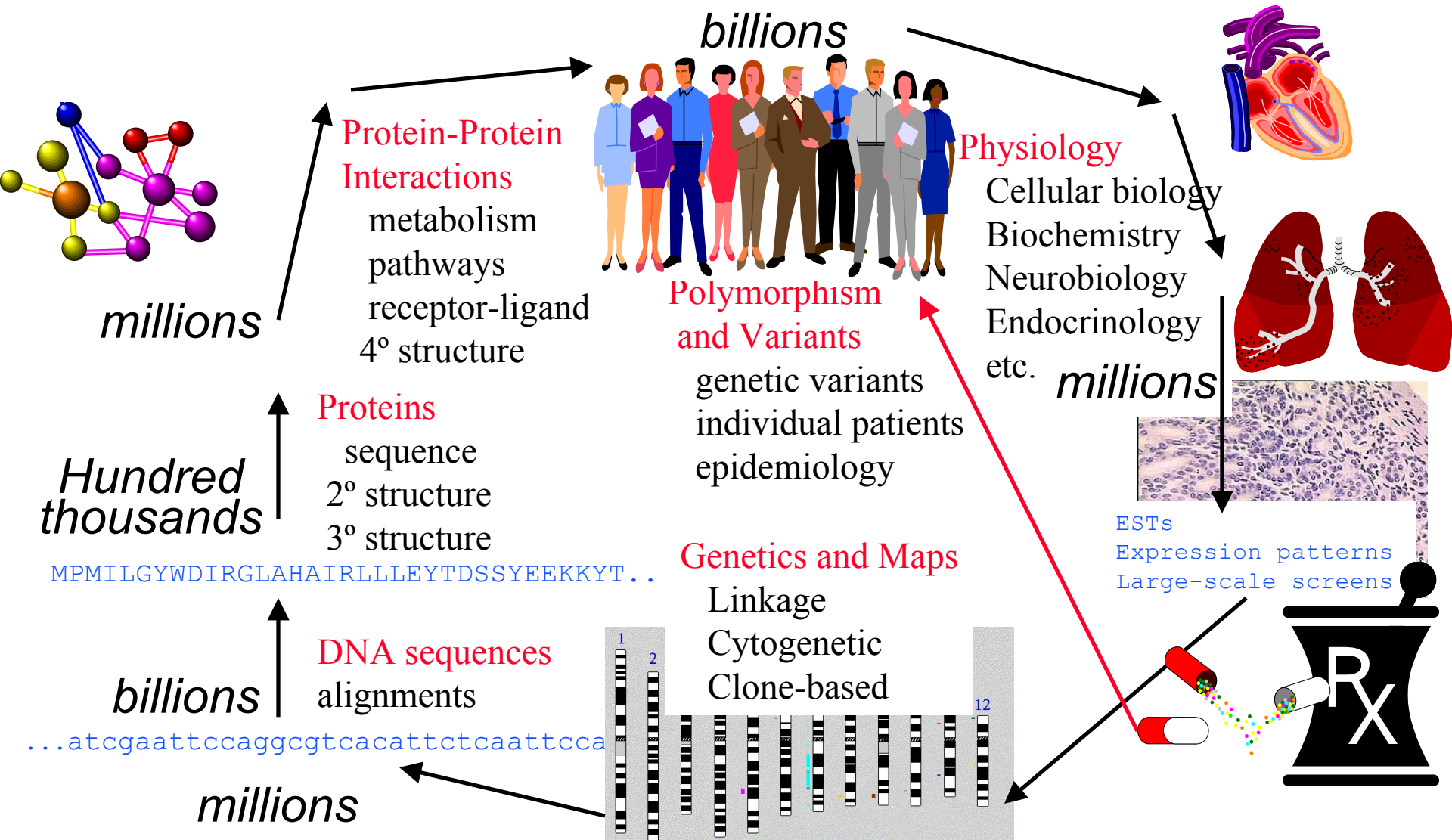
Foreword by STEPHEN JAY GOULD

- Overview of Life on Earth
- Basic Microbiology Introduction
- Excellent Summary

The phyla of life on Earth based on our modification of the Whittaker five-kingdom system and the symbiotic theory of the origin of eukaryotic cells.

PLANTAE
(haplo-diploids)

ANTHOPHYTA

CONIFEROPHYTA

GINKGOPHYTA

GNETOPHYTA

CYCADOPHYTA

FILICINOPHYTA
SPHENOPHYTA
PSILOPHYTA

LYCOPHYTA

BRYOPHYTA

HEPATOPHYTA
ANTHOCEROPHYTA

Histogenesis

Embryogenesis

ANIMALIA
(diploids)

CRANIATA
CEPHALOCHORDATA
UROCHORDATA

ECHINODERMATA
HEMICHORDATA
CHAETOGNATHA

ONYCHOPHORA
TARDIGRADA
MOLLUSCA
POGONOPHORA
ECHIURA
SIPUNCULA
ANNELIDA
CRUSTACEA
MANDIBULATA
CHELICERATA

KINORHYNCHA
ROTIFERA
ACANTHOCEPHALA
NEMATOMORPHA
NEMATODA
NEMERTINA
ORTHONECTIDA
RHOMBOZOA
GNATHOSTOMULIDA
PLATYHELMINTHES

ENTOPROCTA
LORICIFERA
GASTROTRICHA
PRIAPULIDA

CTENOPHORA
CNIDARIA
PLACOZOA

PORIFERA

PHORONIDA
BRACHIOPODA
BRYOZOA

Coelom

FUNGI
(haploids,
dikaryons)

BASIDIOMYCOTA

ASCOMYCOTA

ZYGOMYCOTA

Blastulas

Dikaryosis

PROTOCTISTA

Meiosis

Gametogenesis

RHODOPHYTA
PHAEOPHYTA

CHLOROPHYTA

DISCOMITOCHONDRIA

ZOOMASTIGOTA

CRYPTOMONADA

HAPLOSPORA
APICOMPLEXA
PARAMYXA

GRANULORETICULOSA
XENOPHYOPHORA
MYXOSPORA

LABYRINTHULA

OOMYCOTA
GAMOPHYTA

CHYTRIDIOMYCOTA
PLASMODIOPHORA
HYPHOCHYTRIOMYCOTA
DIATOMS

Mycelia

Mitosis

Steroidogenesis

Aerobiosis

Photosynthesis

Spores

Phototaxis

Motility

XANTHOPHYTA

CHRYSOMONADA

HAPTOMONADA

EUSTIGMATOPHYTA

DINOMASTIGOTA
RHIZOPODA

CILIOPHORA

MICROSPORA
ARCHAEPROTISTA

MYXOMYCOTA

ACTINOPODA

CYANOBACTERIA

ACTINOBACTERIA

CHLOROFLEXA

CHLOROBIA        Chemoautotrophic
                 bacteria

DEINOCOCCI

SAPROSPIRAE

PIRELLULAE

PROTEOBACTERIA

SPIROCHAETAE

ENDOSPORA

THERMOTOGAE

APHRAGMABACTERIA

EUBACTERIA

EURYARCHAEOTA
Methanogens and halophils

CRENARCHAEOTA
Thermoacidophils

Fermentation

ARCHAEA (ARCHAEBACTERIA)

Reproduction

Chromosomal

Chromonemal

From Five Kingdoms by Lynn Margulis and Karlene Schwartz

# Biomedical Data:
## High Complexity and Large Scale

*billions*

*millions*

**Protein-Protein Interactions**
metabolism
pathways
receptor-ligand
4º structure

*Hundred thousands*

**Proteins**
sequence
2º structure
3º structure

MPMILGYWDIRGLAHAIRLLLEYTDSSYEEKKYT...

*billions*

**DNA sequences**
alignments

...atcgaattccaggcgtcacattctcaattcca

*millions*

**Polymorphism and Variants**
genetic variants
individual patients
epidemiology

**Genetics and Maps**
Linkage
Cytogenetic
Clone-based

**Physiology**
Cellular biology
Biochemistry
Neurobiology
Endocrinology
etc.

*millions*

ESTs
Expression patterns
Large-scale screens

# Will Biology Dominate the Grid?

- The largest science discipline:
  - The most scientists (Globally ~500,000-1,000,000)
  - The most research funding (Globally ~$50 Billion/year)
  - The most graduate students (>20,000 year)
- Strong couplings to:
  - Medicine and human health
  - Agriculture and food supplies
  - Energy and ecology
  - Future industrial processes (bio-nano)
  - Consumer of other scientific technologies

# A Diverse Bacterial Community



From Five Kingdoms
Lynn Margulis and Karlene Schwartz

- Pocket in the hindgut wall of the Sonoran desert termite *Pterotermes occidentis*
- 10 billion bacteria per milliliter
- Anoxic environment
- ~30 strains are facultative aerobes
- Many/most are unknown

# Human Microbial Ecology: Sorting Out Cause, Effect + Cure

*Helicobacter pylori* — Stomach cancer
*Helicobacter pylori* — Stomach ulcers
**Cytomegalovirus** — Coronary artery disease
*Porphyromonas gingivalis* — Coronary artery disease
*Chlamydia pneumoniae* — Alzheimer's
**Human papilloma virus** — Cervical cancer
**Hepatitis B, C virus** — Liver cancer
**EpsteinBarr Virus** — Breast cancer
**EpsteinBarr Virus** — Nasopharyngeal carcinoma

— = suspected linkages

Rick Stevens

# An Initial Focus on Prokaryotic life

- ~5,000 known species of prokaryotic life
  - It is estimated that we have identified only 1%-5% of extant species $\Rightarrow$ 80,000-400,000 species
- More diversity than Eukaryotic life forms
  - More diverse metabolisms, more diverse environments
- A Human contains $10^{12}$ cells and $10^{13}$ bacterial cells
  - We have little current understanding of human micro ecology

# Biological CAD: Enabling BioDesign

- Understanding and manipulating biological systems from an information systems standpoint [e.g. organization, communication, transformation]
  - Genotype + Environment = Phenotype
  - Strong analogs to VLSI design tools
- Ultimately goal is designing new biological structures and systems
  - Custom microbe design
  - Metabolic Engineering
  - Synthetic model organisms

Rick Stevens

Argonne ✶ Chicago

*Synechocystis*

3356 Genes
925 Pathways

0.5μm

# Degree of whole-genome structure & function assignment varies by organism (and required confidence level)

## Results for three organisms using Genequiz (EBI)

*Mycoplasma Genitalium*
(468 genes)

*Synechocystis sp.*
(3168 genes)

*Saccharomyces cer.*
(6284 genes)



■ (red) High-confidence structure & function by homology

■ (yellow) Function assigned by strong homology

■ (green) Function assigned by weak homology

■ (cyan) Homologue exists, but function unknown

■ (purple) No homologues

# Reconstructing and Modeling a Prokaryotic Cell?

- ## Typical bacteria [*E. coli*]
  - 1000nm x 300nm x 300nm volume
  - ~4000 genes and gene products
    - 1/4 genes $\Rightarrow$ protein synthesis
    - 1/4 genes $\Rightarrow$ glycolysis
    - 1/4 genes $\Rightarrow$ citric acid cycle
    - rest genes involved in regulation, synthesis and degrading tasks
    - Relatively few genes related to sensing and motility
  - 1,000's of small molecular species, not tracked individually
  - ~3 million total large molecules to track

Rick Stevens

Argonne ✶ Chicago

# Prokaryotic v Eukaryotic Cell



- Membrane separated nucleus
- Multiple cellular compartments
- Introns
- More complex cell membrane
- Complex motility mechanisms
- Etc.

From Lynn Margulis and Karlene Schwartz

LYNN MARGULIS    DORION SAGAN

MICROCOSMOS
FOUR BILLION YEARS OF MICROBIAL EVOLUTION

Foreword by Lewis Thomas

- Lynn Margulis and Dorion Sagan
- Story of microbial evolution
- Introduction to early life on earth

# Intracellular Environment – Gel Like Media



Figure 4.2    Cytoplasm

From: David Goodsell, The Machinery of Life

- 100 nm$^3$
- 450 proteins
- 30 ribosomes
- 340 tRNA molecules
- Several long mRNAs
- 30,000 small organic molecules
- 50,000 Ions
- Rest filled with water 70%

# Cell Membranes and Cell Wall



Figure 4.3    Cell Wall

From: David Goodsell, The Machinery of Life

- Cell wall
  - Polysaccharides
  - Porin pores
- Peptidoglycan
  - cross linked
- Periplasmic space
  - Small proteins
- Complex inner membrane
  - < 50% lipids

# Flagellum and Flagellar Motor: Nanotechnology



Figure 4.4   Flagellum and Flagellar Motor

From: David Goodsell, The Machinery of Life

- Transmembrane proton powered rotating motor
- About 10 Flagella per cell
- 5-10 um long
  - Built from the inside out
- Propels cell ~10-20 um/sec
  - Medium is extremely viscous
  - 10-20 body lengths/sec
  - ~100KM/hr scaled velocity

# DNA Replication via DNA Polymerase



Figure 4.5   Nuclear Region

- <u>DNA replication about 800 new nucleotides per second</u>
- In circular DNA both directions at once
- 50 minute to duplicate entire circle of 4,700,000 nucleotides
- With cell replication ~30 minutes
- <u>DNA replication is pipelined!!</u>

From: David Goodsell, The Machinery of Life

Rick Stevens

Argonne ✶ Chicago

# Understanding Bacterial Life Cycles



Spore formation · Myxospores · Spore germination · Sporangiole · Growing cells · Myxospore formation · Swarming colony · Reproductive structure · Stalk · Germinating cyst

**F** Life cycle of *Stigmatella aurantiaca*. [Drawing by L. Meszoly; labeled by M. Dworkin.]
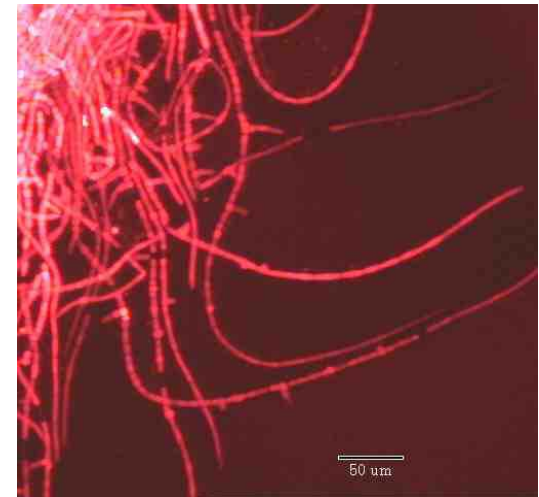
# Modeling Swarming Behavior in *Myxobacteria*



- 100,000 cells swarm to form fruiting bodies
- 80% of the cells lyse
- 20% form spores
- Involves chemotaxis and quorum sensing
- Most complex bacterial genome currently known at > 9Mbp
- Very little is understood

# MONERA: Hierarchical Biological System Modeling Environment

- Genetic Sequences
- Molecular Machines
- Molecular Complexes and modules
- Networks + Pathways [metabolic, signaling, regulation]
- Structural components [ultrastructures]
- Cell Structure and Morphology
- Extracellular Environment
- Populations and Consortia etc.



50 um

Argonne ✶ Chicago

# Systems Biology Model Development and Sharing

| Simulators | Director | Institution | Features |
|---|---|---|---|
| ERATO,j | John Doyle | Caltech | planned workbench |
| Gepasi,w | Pedro Mendes | Santa Fe | MCA, systems kinetics |
| JarnacScamp,wx | Herbert Sauro | Caltech | MCA, Stochastic |
| StochSim,w+ | Dennis Bray | Cambridge | Stochastic |
| BioSpice,u | Adam Arkin | LBL | Stochastic |
| DBSolve,w | Igor Goryanin | Glaxo | enzyme/receptor-ligand |
| E-Cell,u+ | Masaru Tomita | Keio | metabolism. Net ODE |
| Vcell,j | Jim Schaff | U.CT | geometry |
| Xsim,u | J.Bassingthwaighte | Seattle | enzymes to body physiology |
| CellML,x+ | Peter Hunter | U.Auckland | geometry, model sharing |
| GENESIS,u | James Bower | Caltech | neural networks |
| Simex,u+ | Lael Gatewood | U.MN | Stochastic micro populations |
| MONERA, ux+ | Selkov/Stevens | ANL/UC | multilevel ODE/Logic models |

j=java, w=windows, u=unix, x=XML,  + = source/community input  Argonne ✶ Chicago

# What the BioGrid Needs To Provide?

- Scalable compute and data capabilities beyond that available locally

- Distributed infrastructure available 24x7 worldwide

- Integration with local bioinfo systems for seamless computing and data management

- Enables leverage of remote systems administration and support via service providers

- Enables access to state of the art facilities at fraction of the cost (SPs just add more servers)

- Centralized support of tools and data

- Bottom line $\Rightarrow$ <u>enables biologists to focus on biology</u>

# Biology Databases (335 in 2001)

- Major Seq. Repositories (7)
- Comparative Genomics (7)
- Gene Expression (19)
- Gene ID & Structure (31)
- Genetic & Physical Maps (9)
- Genomic (49)
- Intermolecular Interactions (5)
- Metabolic Pathways & Cellular Regulation (12)
- Mutation (34)

- Pathology (8)
- Protein (51)
- Protein Sequence Motifs (18)
- Proteome Resources (8)
- Retrieval Systems & DB Structure (3)
- RNA Sequences (26)
- Structure (32)
- Transgenics (2)
- Varied Biomedical (18)

Baxevanis, A.D. 2002. *Nucleic Acids Research* 30: 1-12.
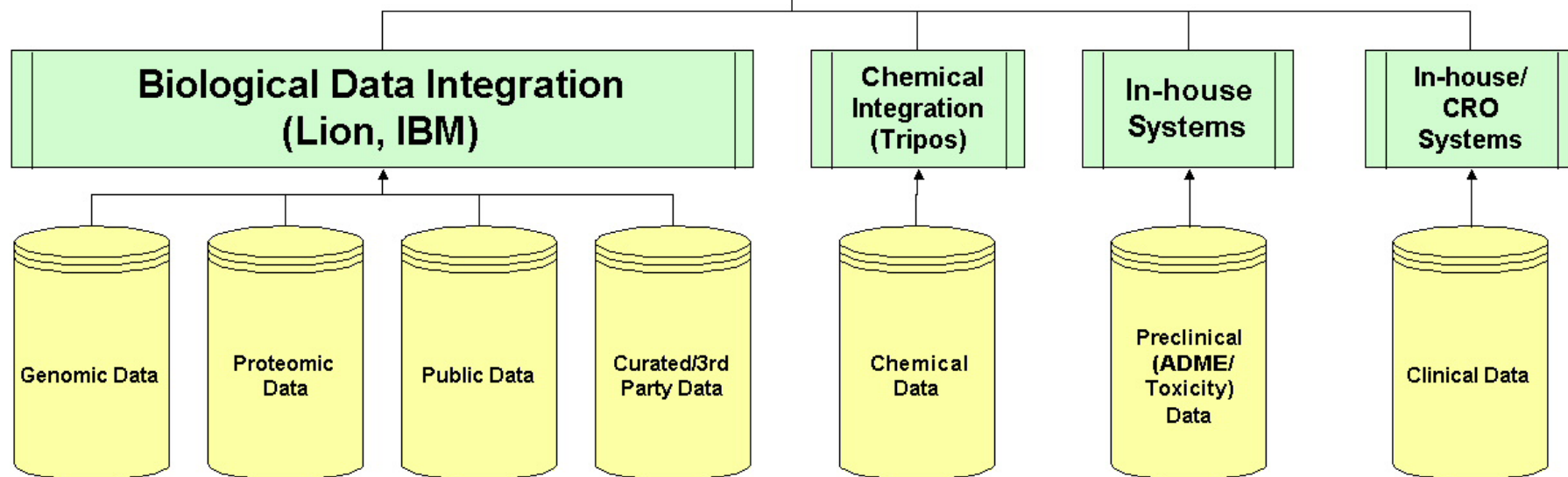
# Grids vs Web tools for biology

- The biology community has developed an extensive collection of web resources to support research:
    - Databases and search engines (entrez, etc)
    - Functional annotation systems (wit, etc.)
    - Organism specific databases (ecocyc, etc.)
    - Literature search engines (pubmed, etc.)
    - Web based modeling systems (vcell, etc.)
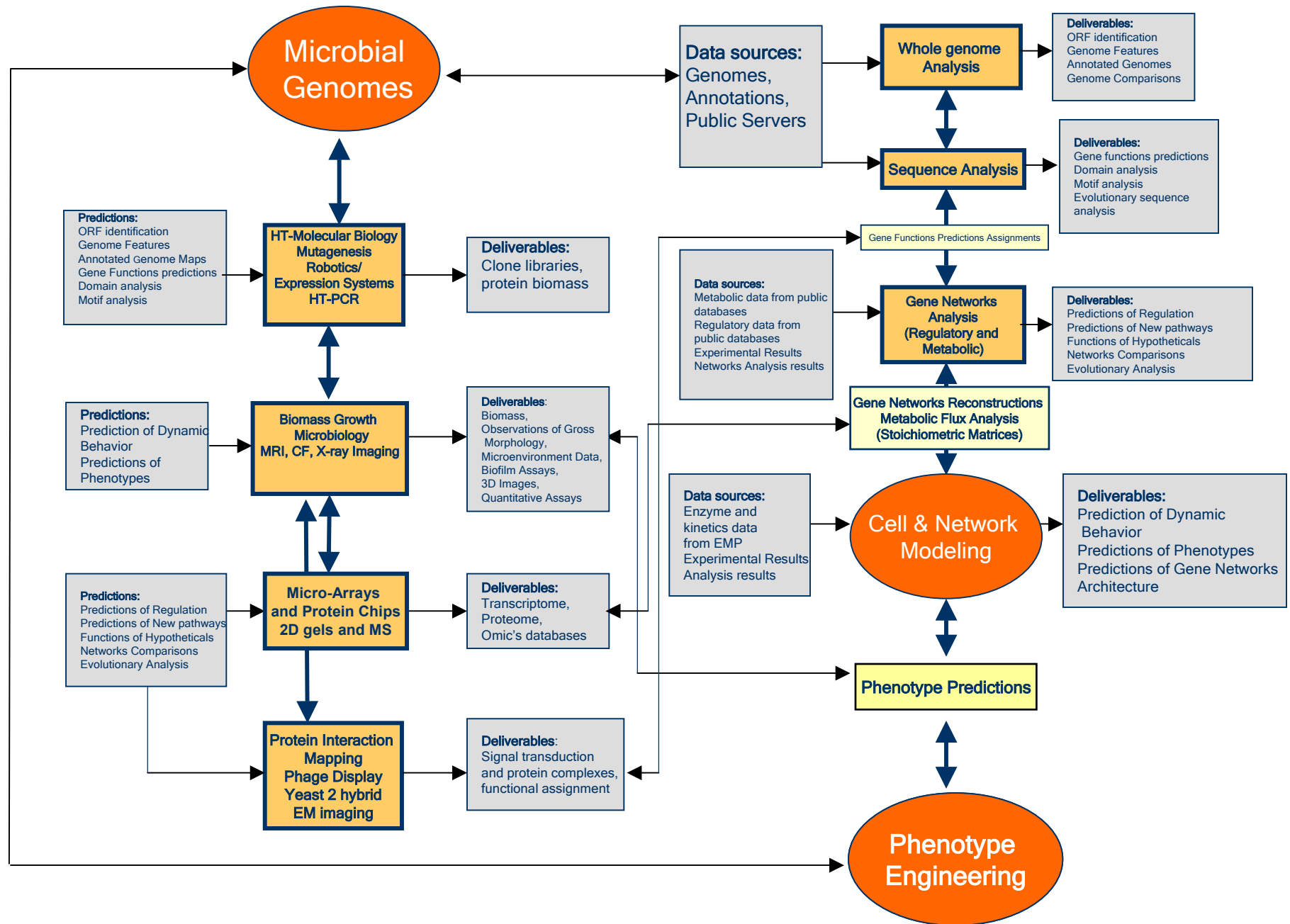
# Software Infrastructure in Drug Discovery



From Richard Gardner(InCellico)
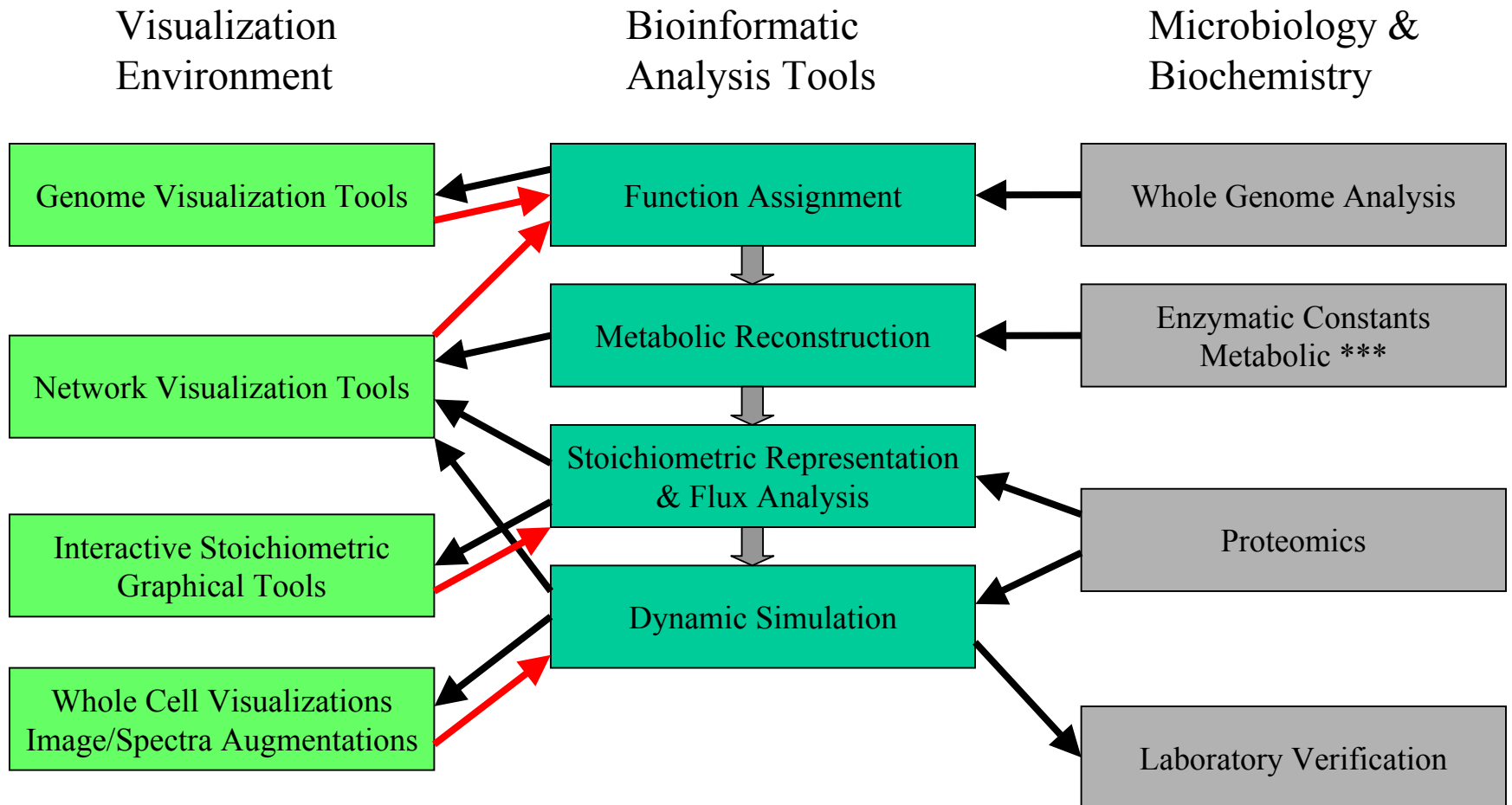
# Requirements for the BioGrid

- Open and extendable architecture
  - Enable tie in to service stack at appropriate points
  - Not just access via Portals
- Leverage scripting tools in wide use for Bioinformatics
  - Create BioGrid services bindings for PERL and Python
- Address data federation and integration
  - Leverage work of IBM, Lion, etc.
- Match the biology workflow and tool chain
  - Create high-level BioGrid services to address critical stages in existing workflow
  - Support composibility of new BioGrid tools with existing tool chain elements

# Visualization + Bioinformatics



Visualization Environment

- Genome Visualization Tools
- Network Visualization Tools
- Interactive Stoichiometric Graphical Tools
- Whole Cell Visualizations Image/Spectra Augmentations

Bioinformatic Analysis Tools

- Function Assignment
- Metabolic Reconstruction
- Stoichiometric Representation & Flux Analysis
- Dynamic Simulation

Microbiology & Biochemistry

- Whole Genome Analysis
- Enzymatic Constants Metabolic ***
- Proteomics
- Laboratory Verification

Rick Stevens

Argonne ✶ Chicago

# Some BioGrid Challenges

- Scalable human bioinformatics expertise
    - Best people working on the important problems
    - Exploit collaboration technology to create world class teams
- Robust local bioinformatics computing environment
    - Best systems administrators and high-end technologies
    - Embed local resources into the Grid via portal technologies
- Access to leading edge bioinformatics software and databases customized to user needs
    - Core content from top scientists and developers
    - Integrated access to biological databases
- Worldwide access to robust computing and database infrastructure
    - Leverage Grid technology to provide worldwide access
    - Integrate purpose built systems and service providers

# What We Need to Create

- Grid Bio applications enablement software layer
  - Provide application's access to Grid services
  - Provides OS independent services
- Grid enabled version of bioinformatics data management tools (e.g. DL, SRS, etc.)
  - Need to support virtual databases via Grid services
  - Grid support for commercial databases
- Bioinformatics applications "plug-in" modules
  - End user tools for a variety of domains
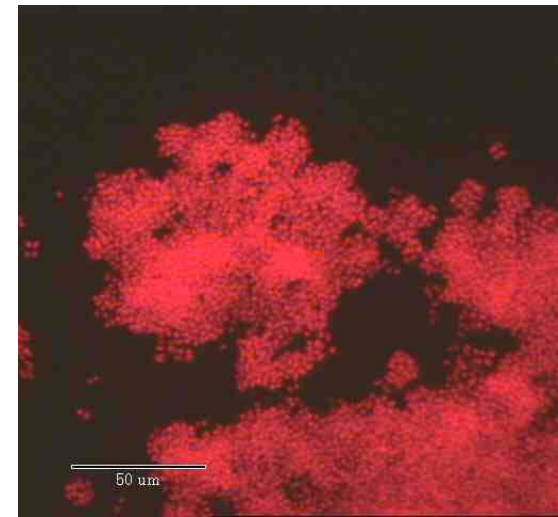  - Support major existing Bio IT platforms

# A Modest Proposal

- Build an Bioinformatics applications development layer on top of basic grid services
  - Think Grid enabled Matlab Toolkit for Biology
- Re-engineer bioinformatics database integration layer to target Grid services model for access
  - Virtualize access to biology databases
- Deploy a network of virtual bioinformatics "Computer Centers" leveraging existing BioGrid resources and new Grid infrastructure (e.g. TeraGrid etc.)
  - Create rich market of resources and services based on common view of BioGrid

# Mathematical Toolkits for Modeling Biological Systems

- "A Mathematica for molecular, cellular and systems biology"
  - Core data models and structures
  - Optimized functions
  - Scripting environment [e.g. Python, PERL, ruby, etc.]
  - Database accessors and built-in schemas
  - Simulation interfaces
  - Parallel and accelerated kernels
  - Visualization interfaces [info-vis and sci-vis]
  - Collaborative workflow and group
    use interfaces



Rick Stevens

# BioGrid Services Model

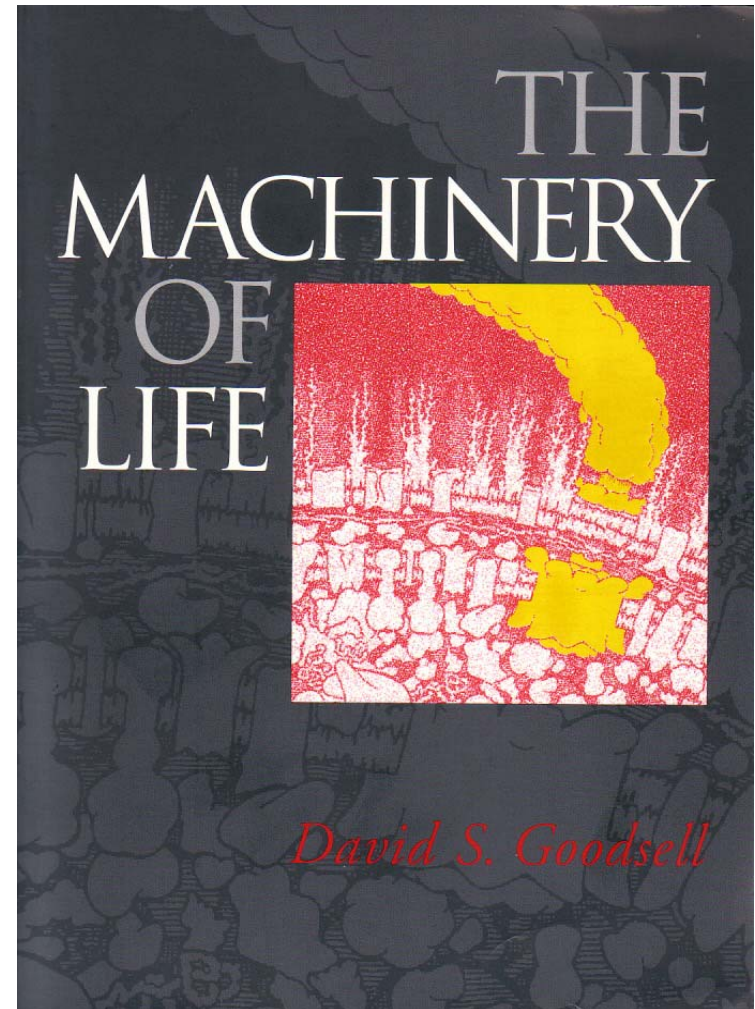| Domain Oriented Services | • Drug Discovery<br>• Microbial Engineering<br>• Molecular Ecology<br>• Oncology Research |
|---|---|
| Basic BioGrid Services | • Integrated Databases<br>• Sequence Analysis<br>• Protein Interactions<br>• Cell Simulation |
| Grid Resource Services | • Compute Services<br>• Pipeline Services<br>• Data Archive Service<br>• Database Hosting |

# An International Systems Biology Grid

- A Data, Experiment and Simulation Grid Linking:
  - People [biologists, computer scientists, mathematicians, etc.]
  - Experimental systems [arrays, detectors, MS, MRI, EM, etc.]
  - Databases [data centers, curators, analysis servers]
  - Simulation Resources [supercomputers, visualization, desktops]
  - Discovery Resources [optimized search servers]
  - Education and Teaching Resources [classrooms, labs, etc.]
- Different than and more fine grain than current Grid Projects
  - More laboratory integration [small laboratory interfaces]
  - Many participants will be experimentalists [workflow, visualization]
  - More diversity of data sources and databases [integration, federation]
  - More portals to simulation environments [ASP models]

- David Goodsell
- Illustrations at 1Mx scale of many cellular domains
- Powerful use of scientific illustration
- Introduction to cellular biology and basic molecular biology

# Conclusions

# Acknowledgements

- DOE, NSF, ANL, UC and Microsoft support my work

- John Wooley (UCSD), Mike Colvin(LLNL/DOE), Ian Foster (ANL/UC), Jack da Silva(NCSC), Bruno Sobral(VT/VBI), Richard Gardner(InCellico)  and others contributed to this talk

# Thank You for Listening

# Science for the 21st Century

- Relevant 20th century milestones

  - Electronic structure theory

  - Silicon based computers

  - Optical networking

  - Software engineering and open source

  - Molecular biology and genomics

  - Electron microscopy and x-ray crystallography

  - Grid computing

# Science for the 21st Century [II]

- Future science milestones
  - First synthetic model prokaryotic organism
  - Characterization of human microbial ecology
  - Global index to life on earth
  - Characterization of microbial life
  - Theory of cell evolution and organization
  - Theory of evolution of intelligence
  - First synthetic eukaryotic organism
  - Confirmation of extra-solar earthlike planets
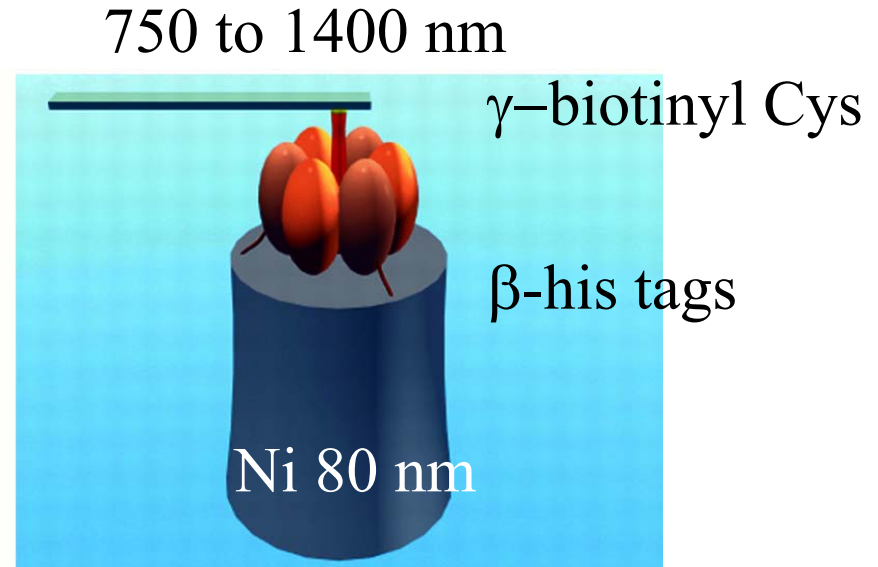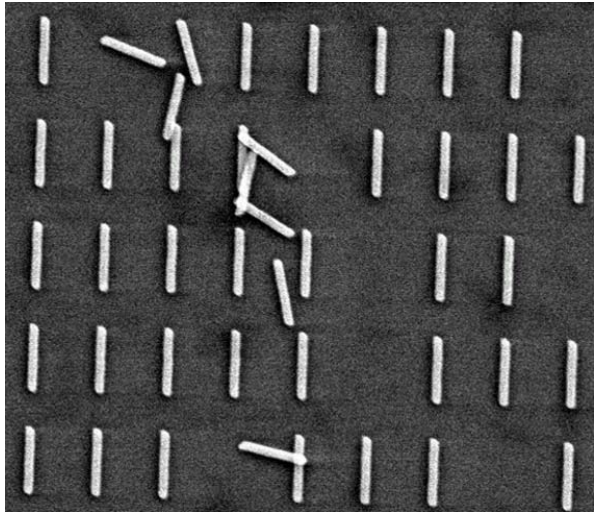  - Synthetic self-reproducing biomemetic nanosystem

# Science for the 21st Century [III]

- The application of advanced biological thought and related technology could yield:
  - Safe and abundant food supplies
  - Sustainable and benign energy sources
  - Effective management of disease and aging
  - Novel materials and renewable industrial feedstocks
  - Advanced computational devices beyond Moore's law
  - Wide variety of molecular scale machinery
  - Self-assembly and self-reproduction technologies
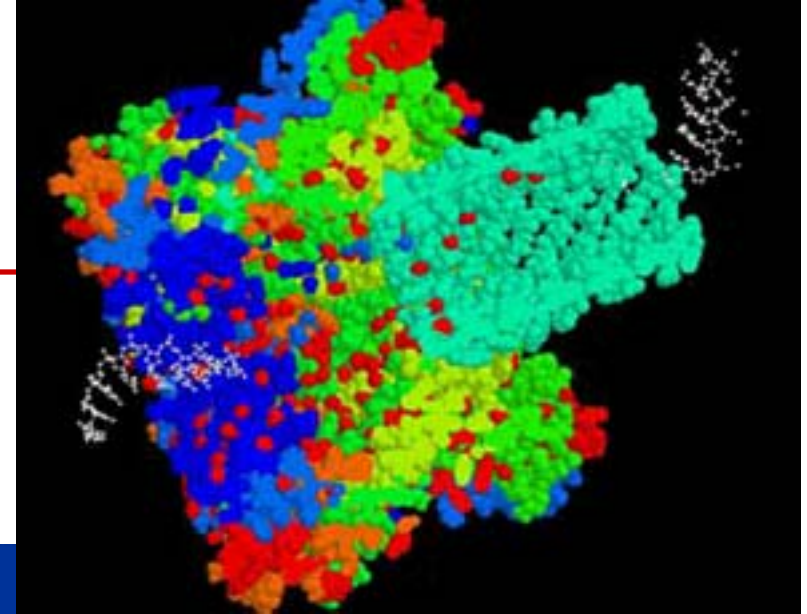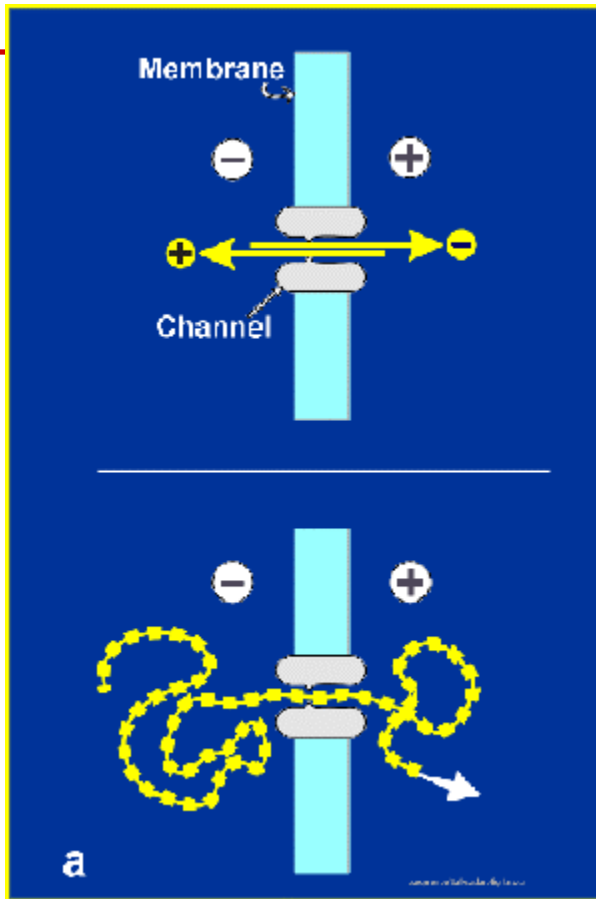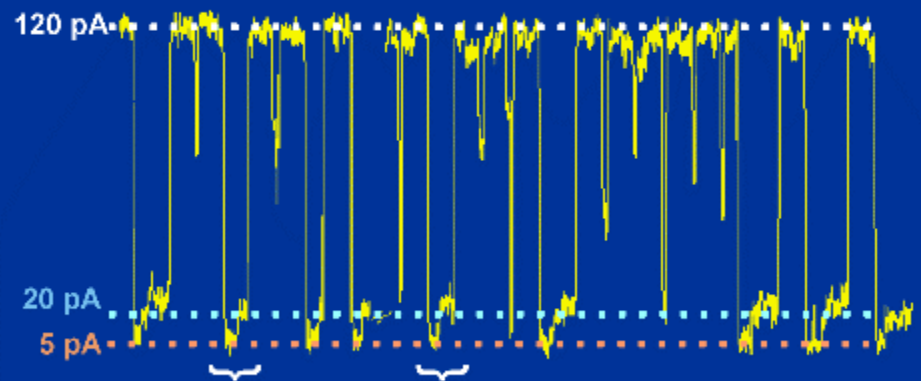
# Nano-ElectroMechanical Systems (NEMS)



750 to 1400 nm

γ−biotinyl Cys

β-his tags

Ni 80 nm

Soong et al. Science 2000; 290: 1555-1558.Powering an Inorganic Nanodevice with a Biomolecular Motor. (Pub)

Rick Stevens

Argonne ✶ Chicago

# Nanosensors
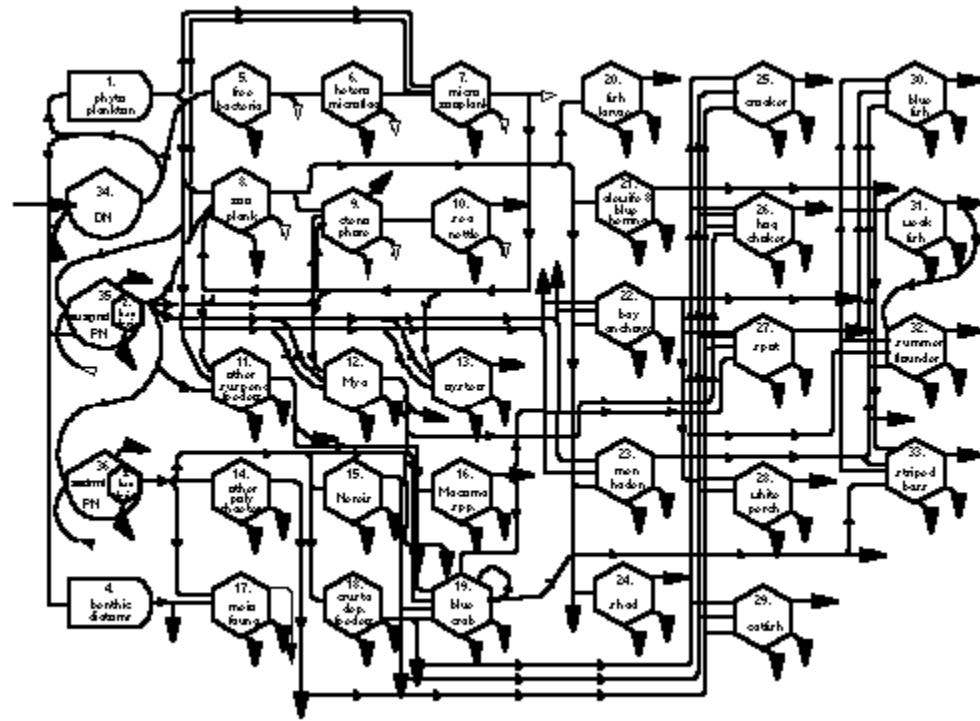




Meller, et al. (2000) "Rapid nanopore discrimination between single polynucleotide molecules." PNAS 1079-84. Akeson et al. Microsecond time-scale discrimination among polyC, polyA, and polyU as homopolymers or as segments within single RNA molecules. Biophys J 1999;77:3227-33

Rick Stevens                                                                        Argonne ✶ Chicago

# Agent Based Approaches to Computational Ecology

- Large-scale ecosystems models
- Individual-base models
  - Age specific behaviors
  - Goal specific behaviors
- Integration with Geochemical process cycles
- Ecosystem Network Analysis
- Microbial ecosystems

# What Could the BioGrid Look Like?

- Logical access to all biological databases
  - Integrated (synthesized?) views of the data with common semantics
  - Variety of active data services
- Transparent access to analysis and modeling services
  - Over 50 commonly used tools with another 100 less common tools
  - Improved composibility of the tools in the tool chain
  - Support for canned analysis and modeling protocols
- Access to a variety of compute and storage resources
  - Services interfaces for computing platforms
  - Grid templates
- Grid enabled scripting languages
  - PERL and Python